

Efficient Estimation of Multidimensional Regression Model using Multilayer Perceptrons

Joseph Rynkiewicz ¹

Université Paris I - SAMOS/MATISSE
90 rue de Tolbiac, Paris - France

February 21, 2008

Abstract

This work concerns the estimation of multidimensional nonlinear regression models using multilayer perceptrons (MLPs). The main problem with such models is that we need to know the covariance matrix of the noise to get an optimal estimator. However, we show in this paper that if we choose as the cost function the logarithm of the determinant of the empirical error covariance matrix, then we get an asymptotically optimal estimator. Moreover, under suitable assumptions, we show that this cost function leads to a very simple asymptotic law for testing the number of parameters of an identifiable MLP. Numerical experiments confirm the theoretical results.

keywords non-linear regression, multivariate regression, multilayer Perceptrons, asymptotic normality

1 Introduction

Let us consider a sequence $(Y_t, Z_t)_{t \in \mathbb{N}}$ of i.i.d. (i.e. independent, identically distributed) random vectors, with Y_t a d -dimensional vector. Each couple (Y_t, Z_t) has the same law as a generic variable (Y, Z) , but it is not hard to generalize all that we show in this paper for stationary mixing variables and therefore for time series. We assume that the model can be written as

$$Y_t = F_{W^0}(Z_t) + \varepsilon_t$$

where

- F_{W^0} is a function represented by an MLP with parameters or weights W^0 .
- (ε_t) is an i.i.d.-centered noise with unknown invertible covariance matrix Γ_0 .

This corresponds to multivariate non-linear least square model, as in chapters 3.1 and 5.1 of Gallant [5]. Indeed, an MLP function can be seen as a parametric non-linear function, for example an one hidden layer MLP using hyperbolic tangent as transfert functions (tanh) can be written $F_{W^0}(Z_t) = (F_{W^0}^1(Z_t), \dots, F_{W^0}^d(Z_t))^T$, where T denotes the transposition of the matrix, with :

$$F_{W^0}^i(z) = \sum_{j=1}^H a_{ij} \tanh \left(\sum_{k=1}^L w_{jk} z_k + w_{j0} \right) + a_{i0}$$

where H is the number of hidden units and L is the dimension of the input z , then the parameter vector is

$$(a_{10}, \dots, a_{dH}, w_{10}, \dots, w_{HL}) \in \mathbb{R}^{(H+1) \times d + (L+1) \times H}$$

There are some obvious transformations that can be applied to an MLP without changing its input-output map. For instance, suppose we pick an hidden node j and we change the sign of all the weights w_{ij} for $i = 0, \dots, H$, and also the sign of all a_{ij} for $i = 0, \dots, d$. Since tanh is odd, this will not alter the contribution of this node to the total net output. Another possibility is to interchange two hidden nodes, that is, to take two hidden nodes j_1 and j_2 and relabel j_1 as j_2 and j_2 as j_1 , taking care to also relabel the corresponding weights. These transformations form a finite group (see Sussmann [10]).

We will consider equivalence classes of one hidden layer MLPs: two MLPs are in the same class if the first one is the image by such transformation of the second one, the considered set of parameters is then the quotient space of parameters by this finite group. In this space, we assume that the model is identifiable it means that the true model belongs to the considered family of models and that we consider MLPs without redundant units. This is a very strong assumption but it is known that estimated weights of an MLPs with redundant units can have a very strange asymptotic behavior (see Kukumizu [4]), because the Hessian matrix is singular. The consequence of the identifiability of the model is that the Hessian matrix computed in the sequel will be definite positive (see Fukumizu [3]). In the sequel we will always assume that we are under the assumptions making the Hessian matrix definite positive.

1.1 Efficient estimation

A popular choice for the associated cost function is the mean square error:

$$\frac{1}{n} \sum_{t=1}^n \|Y_t - F_W(Z_t)\|^2 \quad (1)$$

where $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^d . Although this function is widely used, it is easy to show that we then get a suboptimal estimator, with a larger asymptotic variance than the estimator minimizing the generalized mean square

error :

$$\frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T \Gamma_0^{-1} (Y_t - F_W(Z_t)) \quad (2)$$

But, we need to know the true covariance matrix of the noise to use this cost function. A possible solution is to use an approximation Γ of the covariance error matrix Γ_0 to compute the generalized least squares estimator :

$$\frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T \Gamma^{-1} (Y_t - F_W(Z_t)) \quad (3)$$

A way to construct a sequence of $(\Gamma_k)_{k \in \mathbb{N}^*}$ yielding a good approximation of Γ_0 is the following: using the ordinary least squares estimator \hat{W}_n^1 , the noise covariance can be approximated by

$$\Gamma_1 := \Gamma(\hat{W}_n^1) := \frac{1}{n} \sum_{t=1}^n (Y_t - F_{\hat{W}_n^1}(Z_t))(Y_t - F_{\hat{W}_n^1}(Z_t))^T. \quad (4)$$

then, we can use this new covariance matrix to find a generalized least squares estimator \hat{W}_n^2 :

$$\hat{W}_n^2 = \arg \min_W \frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T (\Gamma_1)^{-1} (Y_t - F_W(Z_t)) \quad (5)$$

and calculate again a new covariance matrix

$$\Gamma_2 := \Gamma(\hat{W}_n^2) = \frac{1}{n} \sum_{t=1}^n (Y_t - F_{\hat{W}_n^2}(Z_t))(Y_t - F_{\hat{W}_n^2}(Z_t))^T.$$

It can be shown that this procedure gives a sequence of parameters

$$\hat{W}_n \rightarrow \Gamma_1 \rightarrow \hat{W}_n^2 \rightarrow \Gamma_2 \rightarrow \dots$$

minimizing the logarithm of the determinant of the empirical covariance matrix (see chapter 5 in Gallant[5]) :

$$U_n(W) := \log \det \left(\frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))(Y_t - F_W(Z_t))^T \right) \quad (6)$$

The use of this cost function for neural networks has been introduced by Williams in 1996 [12], however its theoretical and practical properties have not yet been studied. Here, the calculation of the asymptotic properties of $U_n(W)$ will show that this cost function leads to an asymptotically optimal estimator, with the same asymptotic variance that the estimator minimizing (2), we say then that the estimator is “efficient”.

1.2 testing the number of parameters

Let q be an integer less than s , we want to test “ $H_0 : W \in \Theta_q \subset \mathbb{R}^q$ ” against “ $H_1 : W \in \Theta_s \subset \mathbb{R}^s$ ”, where the sets Θ_q and Θ_s are compact and $\Theta_q \subset \Theta_s$. H_0 expresses the fact that W belongs to a subset Θ_q of Θ_s with a parametric dimension lesser than s or, equivalently, that $s - q$ weights of the MLP in Θ_s are null. If we consider the classical mean square error cost function: $V_n(W) = \sum_{t=1}^n \|Y_t - F_W(Z_t)\|^2$, we get the following test statistic:

$$S_n = n \times \left(\min_{W \in \Theta_q} V_n(W) - \min_{W \in \Theta_s} V_n(W) \right)$$

Under the null hypothesis H_0 , it is shown in Yao [13] that S_n converges in law to a weighted sum of χ_1^2

$$S_n \xrightarrow{\mathcal{D}} \sum_{i=1}^{s-q} \lambda_i \chi_{i,1}^2$$

where the $\chi_{i,1}^2$ are $s - q$ i.i.d. χ_1^2 variables and λ_i are strictly positives eigenvalues of the asymptotic covariance matrix of the estimated weights, different from 1 if the true covariance matrix of the noise is not the identity matrix. So, in the general case, where the true covariance matrix of the noise is not the identity matrix, the asymptotic distribution is not known, because the λ_i s are not known and it is difficult to compute the asymptotic level of the test.

However, if we use the cost function $U_n(W)$ then, under H_0 , the test statistic:

$$T_n = n \times \left(\min_{W \in \Theta_q} U_n(W) - \min_{W \in \Theta_s} U_n(W) \right) \quad (7)$$

will converge to a classical χ_{s-q}^2 so the asymptotic level of the test will be very easy to compute. This is another advantage of using the cost function in Eq. (6). Note that this result is true even if the noise is not Gaussian (it is more general than the maximum likelihood estimator) and without knowing the true covariance of the noise Γ_0 , so without using the cost function (2) or even an approximation of it.

In order to prove these properties, the paper is organized as follows. First we compute the first and second derivatives of $U_n(W)$ with respect to the weights of the MLP, then we deduce the announced properties with classical statistical arguments. Finally, we confirm the theoretical results with numerical experiments.

2 The first and second derivatives of $W \mapsto U_n(W)$

First, we introduce a notation: if $F_W(X)$ is a d -dimensional parametric function depending on a parameter vector W , let us write $\frac{\partial F_W(X)}{\partial W_k}$ (resp. $\frac{\partial^2 F_W(X)}{\partial W_k \partial W_l}$) for the d -dimensional vector of partial derivatives (resp. second order partial

derivatives) of each component of $F_W(X)$. Moreover, if $\Gamma(W)$ is a matrix depending on W , let us write $\frac{\partial}{\partial W_k}\Gamma(W)$ the matrix of partial derivatives of each component of $\Gamma(W)$.

2.1 First derivatives

Now, if $\Gamma_n(W)$ is a matrix depending on the parameter vector W , we get (see Magnus and Neudecker [8])

$$\frac{\partial}{\partial W_k} \log \det (\Gamma_n(W)) = \text{tr} \left(\Gamma_n^{-1}(W) \frac{\partial}{\partial W_k} \Gamma_n(W) \right).$$

Here

$$\Gamma_n(W) = \frac{1}{n} \sum_{t=1}^n (y_t - F_W(z_t))(y_t - F_W(z_t))^T.$$

Note that this matrix $\Gamma_n(W)$ and its inverse are symmetric. Now, if we note that

$$A_n(W_k) = \frac{1}{n} \sum_{t=1}^n \left(-\frac{\partial F_W(z_t)}{\partial W_k} (y_t - F_W(z_t))^T \right),$$

then, using the fact

$$\text{tr} (\Gamma_n^{-1}(W) A_n(W_k)) = \text{tr} (A_n^T(W_k) \Gamma_n^{-1}(W)) = \text{tr} (\Gamma_n^{-1}(W) A_n^T(W_k)),$$

we get

$$\frac{\partial}{\partial W_k} \log \det (\Gamma_n(W)) = 2 \text{tr} (\Gamma_n^{-1}(W) A_n(W_k)). \quad (8)$$

2.2 Calculus of the derivative of $W \mapsto U_n(W)$ for an MLP

Let us note $(\Gamma_n(W))_{ij}$ (resp. $(\Gamma_n^{-1}(W))_{ij}$) the element of the i th line and j th column of the matrix $\Gamma_n(W)$ (resp. $\Gamma_n^{-1}(W)$). We note also $F_W(z_t)(i)$ the i th component of a multidimensional function and for a matrix $A = (A_{ij})$, we note that $(A_{ij})_{1 \leq i, j \leq d}$ is the vector obtained by concatenation of the columns of A . Following the previous results, we can write for the derivative of $\log(\det(\Gamma_n(W)))$ with respect to the weight W_k :

$$\frac{\partial}{\partial W_k} (\log(\det(\Gamma_n(W)))) = \left((\Gamma_n^{-1}(W))_{ij} \right)_{1 \leq i, j \leq d}^T \left(\frac{(\Gamma_n(W))_{ij}}{\partial W_k} \right)_{1 \leq i, j \leq d}$$

with

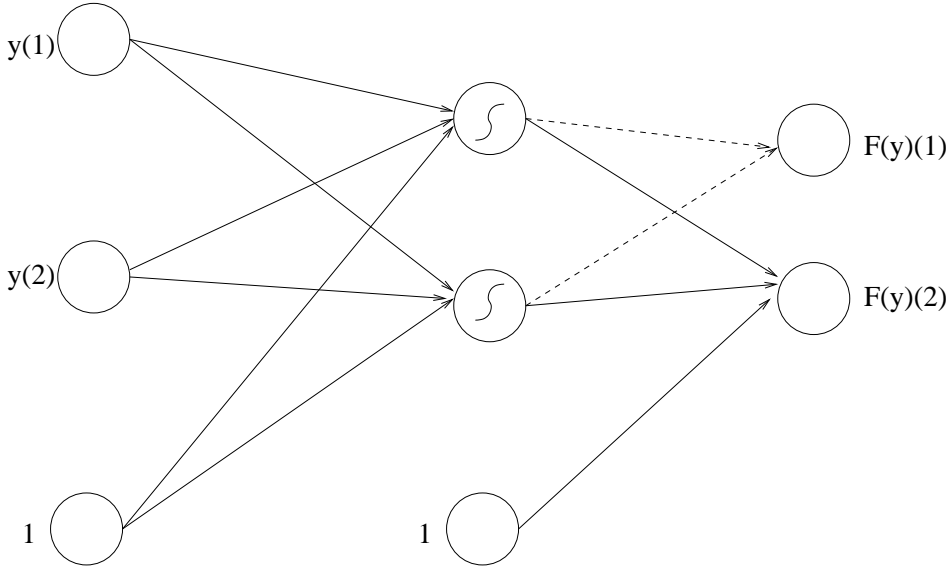
$$\frac{\partial \Gamma_{ij}}{\partial W_k} = \frac{1}{n} \sum_{t=1}^n \left[-\frac{\partial F_W(z_t)(i)}{\partial W_k} \times (y_t - F_W(z_t))(j) - \frac{\partial F_W(z_t)(j)}{\partial W_k} (y_t - F_W(z_t))(i) \right] \quad (9)$$

so

$$\frac{\partial}{\partial W_k}(\log(\det(\Gamma_n(W)))) = \frac{1}{n} (\Gamma_{ij}^{-1})_{1 \leq i, j \leq d}^T \times \left(\sum_{t=1}^n -\frac{\partial F_W(z_t)(i)}{\partial W_k} \times (y_t - F_W(z_t))(j) - \frac{\partial F_W(z_t)(j)}{\partial W_k} (y_t - F_W(z_t))(i) \right)_{1 \leq i, j \leq d} \quad (10)$$

The quantity $\frac{\partial F_W(z_t)(i)}{\partial W_k}$ is computed by back propagating the constant 1 for the MLP restricted to the output i . Figure 1 gives an example of an MLP restricted to the output 2.

Figure 1: MLP restricted to the output 2 : the plain lines



Hence, the calculus of the gradient of $U_n(W)$ with respect to the parameters of the MLP is straightforward. We have to compute the derivative with respect to the weights of each single output MLP extracted from the original MLP by back propagating the constant value 1, then according to the formula (9), we can compute easily the derivative of each term of the empirical covariance matrix of the noise. Finally the gradient is obtained by the sum of all the derivative terms of the empirical covariance matrix multiplied by the terms of its inverse as in formula (10).

2.3 Second derivatives

We write now

$$B_n(W_k, W_l) := \frac{1}{n} \sum_{t=1}^n \left(\frac{\partial F_W(z_t)}{\partial W_k} \frac{\partial F_W(z_t)^T}{\partial W_l} \right)$$

and

$$C_n(W_k, W_l) := \frac{1}{n} \sum_{t=1}^n \left(-(y_t - F_W(z_t)) \frac{\partial^2 F_W(z_t)^T}{\partial W_k \partial W_l} \right)$$

We get

$$\begin{aligned} \frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} &= \frac{\partial}{\partial W_l} 2tr \left(\Gamma_n^{-1}(W) A_n(W_k) \right) = \\ &2tr \left(\frac{\partial \Gamma_n^{-1}(W)}{\partial W_l} A_n(W_k) \right) + 2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n(W)^{-1} C_n(W_k, W_l) \right) \end{aligned}$$

Now, Magnus and Neudecker [8] give an analytic form of the derivative of an inverse matrix, from which we get

$$\frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} = 2tr \left(\Gamma_n^{-1}(W) (A_n(W_k) + A_n^T(W_k)) \Gamma_n^{-1}(W) A_n(W_k) \right) + 2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n^{-1}(W) C_n(W_k, W_l) \right)$$

and

$$\begin{aligned} \frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} &= 4tr \left(\Gamma_n^{-1}(W) A_n(W_k) \Gamma_n^{-1}(W) A_n(W_k) \right) \\ &+ 2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n^{-1}(W) C_n(W_k, W_l) \right) \end{aligned} \quad (11)$$

3 Asymptotic properties

In the sequel, we will assume that the square of the noise ε is integrable and that the cube of the variable Z is integrable too. Moreover, it is easy to show that, for an MLP function, there exists a constant C such that we have the following inequalities :

$$\begin{aligned} \left\| \frac{\partial F_W(Z)}{\partial W_k} \right\| &\leq C(1 + \|Z\|) \\ \left\| \frac{\partial^2 F_W(Z)}{\partial W_k \partial W_l} \right\| &\leq C(1 + \|Z\|^2) \\ \left\| \frac{\partial^2 F_{W^1}(Z)}{\partial W_k \partial W_l} - \frac{\partial^2 F_{W^2}(Z)}{\partial W_k \partial W_l} \right\| &\leq C\|W^1 - W^2\|(1 + \|Z\|^3) \end{aligned}$$

These inequalities will be important to get the local asymptotic normality property implying the asymptotic normality of the parameter minimizing $U_n(W)$.

3.1 Consistency and asymptotic normality of \hat{W}_n

First we have to identify the contrast function associated with $U_n(W)$

Lemma 1

$$U_n(W) - U_n(W^0) \xrightarrow{a.s.} K(W, W^0)$$

with $K(W, W^0) \geq 0$ and $K(W, W^0) = 0$ if and only if $W = W^0$.

Proof: Let us note

$$\Gamma(W) = E((Y - F_W(Z))(Y - F_W(Z))^T)$$

the expectation of the covariance matrix of the noise for model parameter W . By the strong law of large numbers we have

$$U_n(W) - U_n(W^0) \xrightarrow{a.s.} \log \det(\Gamma(W)) - \log \det(\Gamma(W^0)) = \log \frac{\det(\Gamma(W))}{\det(\Gamma(W^0))} = \log \det(\Gamma^{-1}(W^0)(\Gamma(W) - \Gamma(W^0)) + I_d)$$

where I_d denotes the identity matrix of \mathbb{R}^d . So, the lemma is true if $\Gamma(W) - \Gamma(W^0)$ is a positive matrix, null only if $W = W^0$. But this property is true since

$$\begin{aligned} \Gamma(W) &= E((Y - F_W(Z))(Y - F_W(Z))^T) = \\ &= E((Y - F_{W^0}(Z) + F_{W^0}(Z) - F_W(Z))(Y - F_{W^0}(Z) + F_{W^0}(Z) - F_W(Z))^T) = \\ &= E((Y - F_{W^0}(Z))(Y - F_{W^0}(Z))^T) + \\ &= E((F_{W^0}(Z) - F_W(Z))(F_{W^0}(Z) - F_W(Z))^T) = \\ &= \Gamma(W^0) + E((F_{W^0}(Z) - F_W(Z))(F_{W^0}(Z) - F_W(Z))^T) \end{aligned}$$

and the lemma follows from the identifiability assumption ■

We deduce the theorem of consistency:

Theorem 1 *We have*

$$\hat{W}_n \xrightarrow{a.s.} W^0$$

Proof Remark that a constant B exists such that

$$\sup_{W \in \Theta_s} \|Y - F_W(Z)\|^2 < \|Y\|^2 + B \quad (12)$$

because Θ_s is compact, so $F_W(Z)$ is bounded. Let us define the function

$$\Phi(\Gamma) := \max(\log \det(\Gamma), d \log(\delta))$$

where d is the dimension of the observations Y and $\delta > 0$ strictly smaller than the smallest eigenvalue of Γ_0 , since Γ_0 is definite positive we have for all W :

$$\lim_{n \rightarrow \infty} \Phi(\Gamma_n(W)) \stackrel{a.s.}{=} \lim_{n \rightarrow \infty} \log \det(\Gamma_n(W)) = K(W, W^0) + \log \det(\Gamma^0) > d \log(\delta)$$

Now, for all W , thanks to the inequality (12) there exists constants α and β such that

$$|\Phi((Y - F_W(Z))(Y - F_W(Z))^T)| \stackrel{a.s.}{\leq} \alpha \|Y\|^2 + \beta$$

but the right hand of this inequality is integrable, so the function Φ as an integrable envelope function and by example 19.8 of van der Vaart [11] the set of functions $\{\Phi((Y - F_W(Z))(Y - F_W(Z))^T), W \in \Theta_s\}$ is Glivenko-Cantelli.

Now, the theorem 5.7 of van der Vaart [11], shows that \hat{W}_n converges in probability to W^0 , but it is easy to show that this convergence is almost sure.

First, by lemma 1, we remark that for all neighborhood \mathcal{N} of W^0 there exists a number $\eta(\mathcal{N}) > 0$ such that for all $W \notin \mathcal{N}$ we have

$$\log \det (\Gamma(W)) > \log \det (\Gamma(W^0)) + \eta(\mathcal{N})$$

Now to show the strong consistency property we have to prove that for all neighborhood \mathcal{N} of W^0 we have $\lim_{n \rightarrow \infty} \hat{W}_n \stackrel{a.s.}{\subset} \mathcal{N}$ or, equivalently,

$$\lim_{n \rightarrow \infty} \log \det (\Gamma(\hat{W}_n)) - \log \det (\Gamma(W^0)) < \eta(\mathcal{N})$$

By definition, we have

$$\log \det (\Gamma_n(\hat{W}_n)) \leq \log \det (\Gamma_n(W^0))$$

and the Glivenko-Cantelli property assures that

$$\lim_{n \rightarrow \infty} \log \det (\Gamma_n(W^0)) - \log \det (\Gamma(W^0)) \stackrel{a.s.}{=} \lim_{n \rightarrow \infty} \Phi(\Gamma_n(W)) - \log \det (\Gamma(W^0)) \stackrel{a.s.}{=} 0$$

therefore

$$\lim_{n \rightarrow \infty} \log \det (\Gamma_n(\hat{W}_n)) < \log \det (\Gamma(W^0)) + \frac{\eta(\mathcal{N})}{2}$$

We have also

$$\lim_{n \rightarrow \infty} \log \det (\Gamma_n(\hat{W}_n)) - \log \det (\Gamma(\hat{W}_n)) \stackrel{a.s.}{=} \lim_{n \rightarrow \infty} \Phi(\Gamma_n(W)) - \log \det (\Gamma(\hat{W}_n)) \stackrel{a.s.}{=} 0$$

and finally

$$\lim_{n \rightarrow \infty} \log \det (\Gamma(\hat{W}_n)) - \frac{\eta(\mathcal{N})}{2} < \log \det (\Gamma_n(\hat{W}_n)) < \log \det (\Gamma(W^0)) + \frac{\eta(\mathcal{N})}{2}$$

■

Now, we can establish the asymptotic normality for the estimator.

Lemma 2 *Let $\Delta U_n(W^0)$ be the gradient vector of $U_n(W)$ at W^0 , $\Delta U(W^0)$ be the gradient vector of $U(W) := \log \det (\Gamma(W))$ at W^0 and $HU_n(W^0)$ be the Hessian matrix of $U_n(W)$ at W^0 .*

We define finally

$$B(W_k, W_l) := \frac{\partial F_W(Z)}{\partial W_k} \frac{\partial F_W(Z)}{\partial W_l}^T$$

Then we get

1. $HU_n(W^0) \xrightarrow{a.s.} 2I_0$
2. $\sqrt{n}\Delta U_n(W^0) \xrightarrow{Law} \mathcal{N}(0, 4I_0)$

where, the component (k, l) of the matrix I_0 is :

$$tr (\Gamma_0^{-1} E (B(W_k^0, W_l^0)))$$

proof First we note

$$A(W_k) = \left(-\frac{\partial F_W(Z)}{\partial W_k} (Y - F_W(Z))^T \right)$$

To prove the lemma, we remark first that the component (k, l) of the matrix $4I_0$ is :

$$E \left(\frac{\partial U(W^0)}{\partial W_k} \frac{\partial U(W^0)}{\partial W_l^0} \right) = E \left(2tr \left(\Gamma_0^{-1} A^T(W_k^0) \right) \times 2tr \left(\Gamma_0^{-1} A(W_l^0) \right) \right)$$

and, since the trace of the product is invariant by circular permutation,

$$\begin{aligned} E \left(\frac{\partial U(W^0)}{\partial W_k} \frac{\partial U(W^0)}{\partial W_l^0} \right) &= \\ 4E \left(-\frac{\partial F_{W^0}(Z)^T}{\partial W_k} \Gamma_0^{-1} (Y - F_{W^0}(Z)) (Y - F_{W^0}(Z))^T \Gamma_0^{-1} \left(-\frac{\partial F_{W^0}(Z)}{\partial W_l} \right) \right) \\ &= 4E \left(\frac{\partial F_{W^0}(Z)^T}{\partial W_k} \Gamma_0^{-1} \frac{\partial F_{W^0}(Z)}{\partial W_l} \right) \\ &= 4tr \left(\Gamma_0^{-1} E \left(\frac{\partial F_{W^0}(Z)}{\partial W_k} \frac{\partial F_{W^0}(Z)^T}{\partial W_l} \right) \right) \\ &= 4tr \left(\Gamma_0^{-1} E \left(B(W_k^0, W_l^0) \right) \right) \end{aligned}$$

Now, for the component (k, l) of the expectation of the Hessian matrix, we remark that

$$\lim_{n \rightarrow \infty} tr \left(\Gamma_n^{-1}(W^0) A_n(W_k^0) \Gamma_n^{-1}(W^0) A_n(W_k^0) \right) = 0$$

and

$$\lim_{n \rightarrow \infty} tr \Gamma_n^{-1} C_n(W_k^0, W_l^0) = 0$$

so

$$\begin{aligned} \lim_{n \rightarrow \infty} H_n(W^0) &= \lim_{n \rightarrow \infty} 4tr \left(\Gamma_n^{-1}(W^0) A_n(W_k^0) \Gamma_n^{-1}(W^0) A_n(W_k^0) \right) + \\ &2tr \Gamma_n^{-1}(W^0) B_n(W_k^0, W_l^0) + 2tr \Gamma_n^{-1} C_n(W_k^0, W_l^0) = \\ &= 2tr \left(\Gamma_0^{-1} E \left(B(W_k^0, W_l^0) \right) \right) \end{aligned}$$

■

Now, from a classical argument of local asymptotic normality (see for example Yao [13]), we deduce the following property for the estimator \hat{W}_n :

Proposition 1 *We have*

$$\lim_{n \rightarrow \infty} \sqrt{n}(\hat{W}_n - W^0) = \mathcal{N}(0, I_0^{-1})$$

However, if W_n^* is the estimator of the generalized least squares :

$$W_n^* := \arg \min \frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T \Gamma_0^{-1} (Y_t - F_W(Z_t))$$

then we have also

$$\lim_{n \rightarrow \infty} \sqrt{n}(W_n^* - W^0) = \mathcal{N}(0, I_0^{-1})$$

so \hat{W}_n has the same asymptotic behavior as the generalized least squares estimator with the true covariance matrix Γ_0^{-1} which is asymptotically optimal (see for example Ljung [7]). Therefore, the proposed estimator is asymptotically optimal too.

3.2 Asymptotic distribution of the test statistic T_n

Let us assume that the null hypothesis H_0 is true, we write $\hat{W}_n = \arg \min_{W \in \Theta_s} U_n(W)$ and $\hat{W}_n^0 = \arg \min_{W \in \Theta_q} U_n(W)$, where Θ_q is viewed as a subset of Θ_s . The asymptotic distribution of T_n is then a consequence of the previous section. Namely, if we replace $nU_n(W)$ by its Taylor expansion around \hat{W}_n and \hat{W}_n^0 , following van der Vaart [11] chapter 16 we have :

$$T_n = \sqrt{n} \left(\hat{W}_n - \hat{W}_n^0 \right)^T I_0 \sqrt{n} \left(\hat{W}_n - \hat{W}_n^0 \right) + o_P(1) \xrightarrow{\mathcal{D}} \chi_{s-q}^2$$

4 Experimental results

4.1 Simulated example

Although the estimator associated with the cost function $U_n(W)$, is theoretically better than the ordinary mean least squares estimator, it is of some interest to quantify this fact by simulation. Moreover, there are some pitfalls in practical situations with MLPs.

The first point is that we have no guaranty to reach the global minimum of the cost function, we can only hope to find a good local minimum if we are using many estimations with different initial weights.

The second point, is the fact that MLP are black box, it means that it is difficult to give an interpretation of their parameters and it is almost impossible to compare MLP by comparing their parameters even if we try to take into account the possible permutations of the weights, because the difference between the weights may reflect only the differences of local minima reached by weights during the learning.

All these reasons explain why we choose, for simplicity, to compare the estimated covariance matrices of the noise instead of comparing directly the estimated parameters of MLPs.

4.1.1 The model

To simulate our data, we use an MLP with 2 inputs, 3 hidden units, and 2 outputs. We choose to simulate an auto-regressive time series, where the outputs at time t are the inputs for time $t + 1$. Moreover, with MLPs, the statistical properties of such a model are the same as with independent identically distributed (i.i.d.) data, because the time series constitutes a mixing process (see Yao [13]).

The equation of the model is the following

$$Y_{t+1} = F_{W_0}(Y_t) + \varepsilon_{t+1}$$

where

- $Y_0 = (0, 0)$.
- $(Y_t)_{1 \leq t \leq 1000}$, $Y_t \in \mathbb{R}^2$, is the bidimensional simulated random process
- F_{W_0} is an MLP function with weights W_0 chosen randomly between -2 and 2 .
- (ε_t) is an i.i.d. centered noise with covariance matrix $\Gamma_0 = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}$.

In order to study empirically the statistical properties of our estimator we make 400 independent simulations of the bidimensional time series of length 1000.

4.1.2 The results

Our goal is to compare the estimator minizing $U_n(W)$ or equation (6) and the weights minimizing the mean square error (MSE), equation (1). For each time series we estimate the weights of the MLP using the cost function $U_n(W)$ and the MSE. The estimations have been done using the second order algorithm BFGS, and for each estimation we choose the best result obtained after 100 random initializations of the weights. Thus, we avoid plaguing our learning with poor local minima.

We show here the mean of estimated covariance matrices of the noise for $U_n(W)$ and the mean square error (MSE) cost function:

$$U_n(W) : \begin{pmatrix} 4.960 & 3.969 \\ 3.969 & 4.962 \end{pmatrix} \text{ and MSE : } \begin{pmatrix} 4.938 & 3.932 \\ 3.932 & 4.941 \end{pmatrix}$$

The estimated standard deviation of the terms of the matrices are all equal to 0.01, so the differences observed between the two matrices are statistically significant. We can see that the estimated covariance of the noise is on average better with the estimator associated to the cost function $U_n(W)$, in particular it seems that there is slightly less overfitting with this estimator, and the non diagonal terms are greater than with the least squares estimator. As expected, the determinant of the mean matrix associated with $U_n(W)$ is 8.86 instead of 8.93 for the matrix associated with the MSE.

4.2 Application to real time series: Pollution of ozone

Ozone is a reactive oxidant, which is formed both in the stratosphere and troposphere. Near the ground's surface, ozone is directly harmful to human health, plant life and damages physical materials. The population, especially in large cities and in suburban zones which suffer from summer smog, wants to be warned

of high pollutant concentrations in advance. The statistical ozone modelling and more particularly regression models have been widely studied [1], [6]. Generally, linear models do not seem to capture all the complexity of the phenomena. Thus, the use of nonlinear techniques is recommended to deal with ozone prediction. Here we want to predict ozone pollution in two sites at the same time. The sites are the pollution levels in the south of Paris (13th district) and on the top of the Eiffel Tower. As these sites are very near each other we can expect that the two components of the noise are very correlated.

4.2.1 The model

The neural model used in this study is autoregressive and includes exogenous parameters (called NARX model), where X stands for exogeneous variables. Our aim is to predict the maximum level of ozone pollution of the next day knowing the today's maximum level of pollution and the maximal temperature of the next day. If we note Y^1 the level of pollution for Paris 13, Y^2 the level of pollution for the Eiffel Tower and $Temp$ the temperature, the model can be written as follows:

$$(Y_{t+1}^1, Y_{t+1}^2) = F_W(Y_t^1, Y_t^2, Temp_{t+1}) + \varepsilon_{t+1} \quad (13)$$

We will assume that the variables are mixing as previously. As usual with real time series, overtraining is a crucial problem. MLPs are very overparametrized models. This occurs when the model learns the details of the noise of the training data. Overtrained models have very poor performance on fresh data. To avoid overtraining we use in this study the SSM pruning technique, a statistical stepwise method using a BIC-like criterion (Cottrell et al [2]). The MLP with the minimal dimension is found by the elimination of the irrelevant weights. Here, we will compare behavior of this method for both cost function: The mean square error (MSE) and the logarithm of the determinant of the empirical covariance matrix of the noise ($U_n(W)$).

4.2.2 The dataset

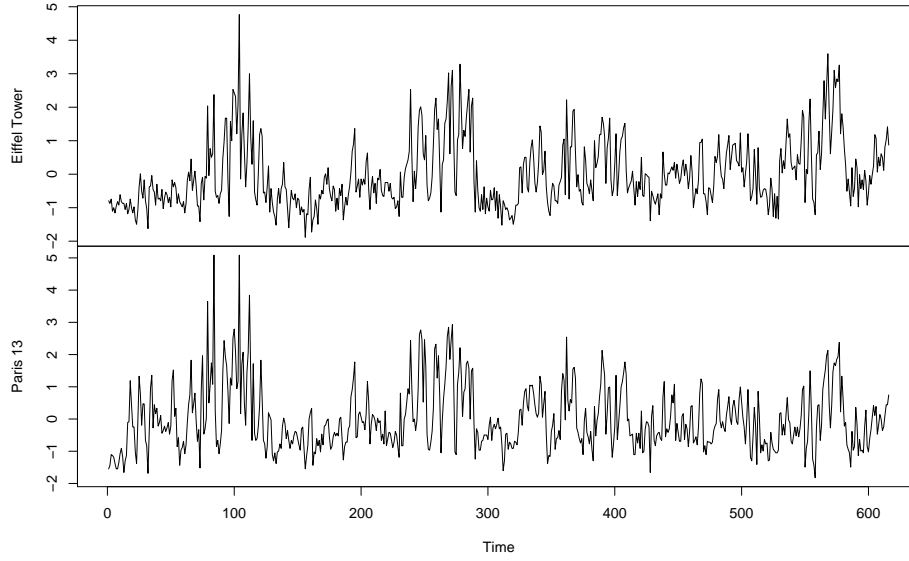
This study presents the ozone concentration of the Air Quality Network of the Ile de France Region (AIRPARIF, Paris, France). The data used in this work are from 1994 to 1997, we use only the months from April to September inclusive because there is no peak during the winter period. According to the model, we have the following parameters:

- The maximum temperature of the day
- Persistence is used by introducing the previous day's peak ozone.

Before their use in the neural network, all these data have been centered and normalized. The data used to train the MLPs are chosen randomly in the whole period and we leave 100 observations to form a fresh data set (test set), which will be used for models evaluation. In order to evaluate the models we repeat

400 times this random sampling to get 400 covariance matrices on each set for the two cost functions. Figure 2 is a plot of the centered and normalized original data.

Figure 2: Ozone time series



4.2.3 The results

For the learning set, we get the following results for the averaged covariance matrix (the estimated standard deviation for the coefficients is about 0.0005):

$$U_n(W) : \begin{pmatrix} 0.27 & 0.20 \\ 0.20 & 0.34 \end{pmatrix} \text{ and } \text{MSE} : \begin{pmatrix} 0.27 & 0.18 \\ 0.18 & 0.34 \end{pmatrix}$$

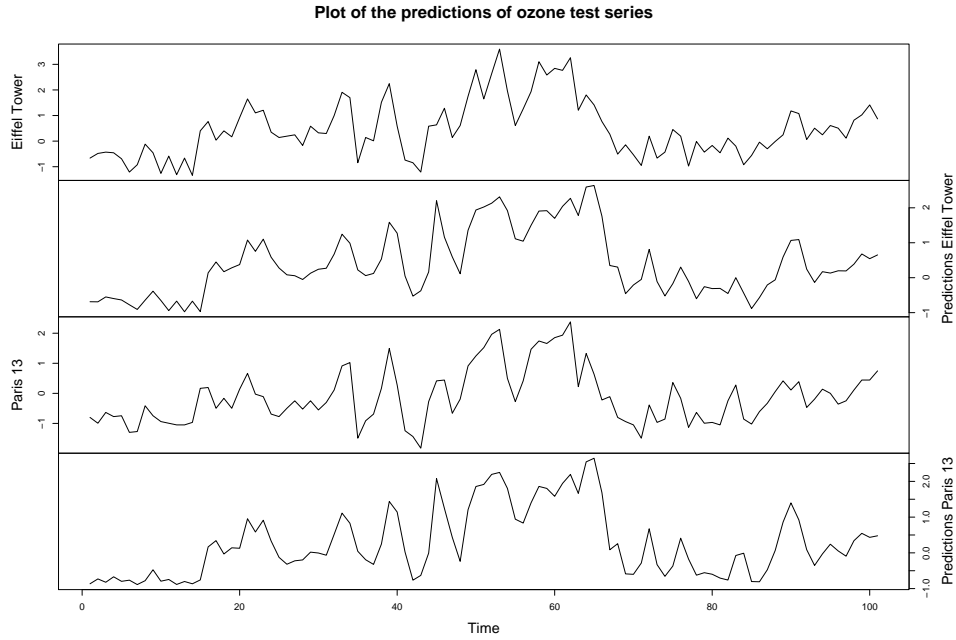
for the test set, we get the following results for the averaged covariance matrix (the estimated standard deviation for the coefficients is about 0.002) :

$$U_n(W) : \begin{pmatrix} 0.29 & 0.22 \\ 0.22 & 0.36 \end{pmatrix} \text{ and } \text{MSE} : \begin{pmatrix} 0.33 & 0.20 \\ 0.20 & 0.39 \end{pmatrix}$$

The two matrices are almost the same for the learning set, however the non-diagonal terms are greater for the $U_n(W)$ cost function. Moreover, looking at

the averaged matrix on the test set, we see that the generalization capabilities are better for $U_n(W)$ and the differences are statistically significant. Generally, the best MLP for $U_n(W)$ has less weights than the best MLP for the MSE cost function. Hence, the proposed cost function leads to a somewhat more parsimonious model, because the pruning technique is very sensitive to the variance of estimated parameters. This gain is valuable regarding the generalization capacity of the model, because the difference is almost null for the learning data set but is greater on the test data. Figure 3 is a plot of the centered and normalized original test data and its prediction.

Figure 3: Predicted time series



5 Conclusion

In the linear multidimensional regression model the optimal estimator has an analytic solution (see Magnus and Neudecker [8]), so it does not make sense to consider minimization of a cost function. However, for the non-linear multidimensional regression model, the ordinary least squares estimator is sub-optimal, if the covariance matrix of the noise is not the identity matrix. We can over-

come this difficulty by using the cost function $U_n(W) = \log \det(\Gamma_n(W))$. In this paper, we have provided a proof of the optimality of the estimator associated with $U_n(W)$. Statistical thought tells us that it is always better for the neural networks practitioners to use a more efficient estimator because such estimator are better on average, even if the difference seems to be small. This estimator is especially important if the practitioners are using pruning techniques. Indeed pruning technique are based on Wald test or approximated Wald test as for the optimal brain damage or optimal brain surgeon method (see Cottrell et al. [2]) and these tests are very sensitive to the variance of the estimated parameters. Moreover, we have shown that this cost function leads to a simpler χ^2 test to determine the number of weights if the model is identifiable. These theoretical results have been confirmed by a simulated example, and we have seen for a real time series that we can expect slight improvement especially in model selection, this confirms the fact that such techniques are very sensitive to the variance of the estimated weights.

References

- [1] A.C. Comrie, Comparing neural networks and regression models for ozone forecasting, Air and Waste Management Association, 47 (1997) 653-663.
- [2] M. Cottrell, et al., Neural modeling for time series: a statistical stepwise method for weight elimination, IEEE Transaction on Neural Networks 6 (1995) 1355-1364.
- [3] K. Fukumizu, A regularity condition of the information matrix of a multilayer perceptron network, Neural Networks, 9:5 (1996) 871-879.
- [4] K. Fukumizu, Likelihood ratio of unidentifiable models and multilayer neural networks, The Annals of Statistics, 31:3 (2003) 833-851.
- [5] R. A. Gallant, Non linear statistical models (J. Wiley and Sons, New York, 1987).
- [6] M.W. Gardner and S.R. Dorling, Artificial neural networks, the multilayer Perceptron. A review of applications in the atmospheric sciences, Atmospheric Environment, 32:14/15 (1998) 2627-2636.
- [7] L. Ljung, System identification: Theory for the user (Prentice Hall, New Jersey, 1999).
- [8] J. Magnus and H. Neudecker, Matrix differential calculus with applications in statistics and econometrics (J. Wiley and Sons, New York, 1988).
- [9] J. Rynkiewicz, Estimation of Multidimensional Regression Model with Multilayer Perceptron, in: J. Mira and A. Prieto, ed., proc. IWANN'2003, Lecture Notes in Computer Science, Vol. 2686 (Springer, Berlin, 2003) 310-317.
- [10] H.J. Sussmann, Uniqueness of the weights for minimal feedforward nets with a given input-output Map, Neural Networks 5 (1992) 589-593.
- [11] A. W. Van der Vaart, Asymptotic statistics (Cambridge University Press, Cambridge, 1998).
- [12] P. M. Williams, Using neural networks to model conditional multivariate densities, Neural Computation 8:4 (1996) 843-854.
- [13] J.F. Yao, On least squares estimation for stable nonlinear AR processes, The Annals of the Institute of Mathematical Statistics 52 (2000) 316-331.